

Sequence Models

Examples of sequence data / tasks :

sequence recognition

Music generation

Sentiment classification

DNA sequence analysis

Machine translation

Video activity recognition

Name Entity Recognition

either x , or y , or both
are sequence data.

e.g. Sequence of words,

sequence of notes ..

Notation

Motivating Example:

x : (Harry Potter) and (Hermione Granger) invented a new spell

$x^{<1>} \quad x^{<2>} \quad \dots$

$x^{<t>} \quad x^{<9>} \quad T_x = 9$

$\rightarrow y$: 1 1 0 1 1 0 0 0 0

$y^{<1>} \quad y^{<2>} \quad \dots$

$y^{<9>} \quad T_y = 9$

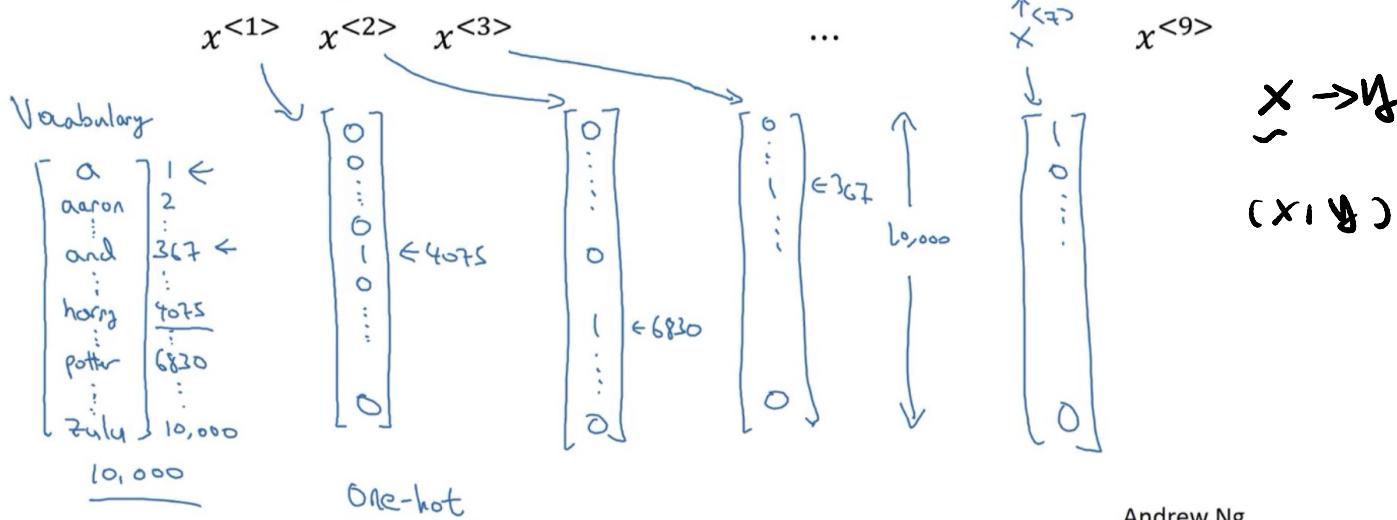
$x^{(i)<t>} : t^{\text{th}}$ element in
 i^{th} training example

$T_x^{(i)}$: input length for training example i

$T_y^{(i)}$

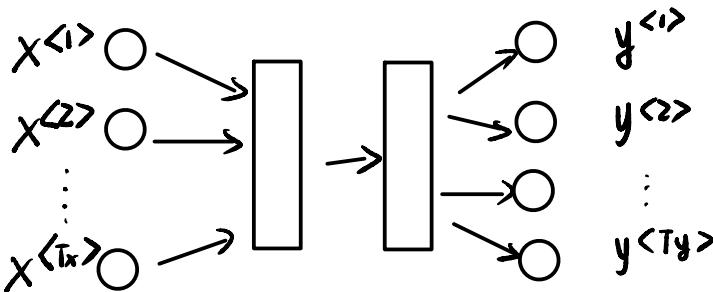
Representing Words

x : Harry Potter and Hermione Granger invented a new spell.



Recurrent Neural Network Model

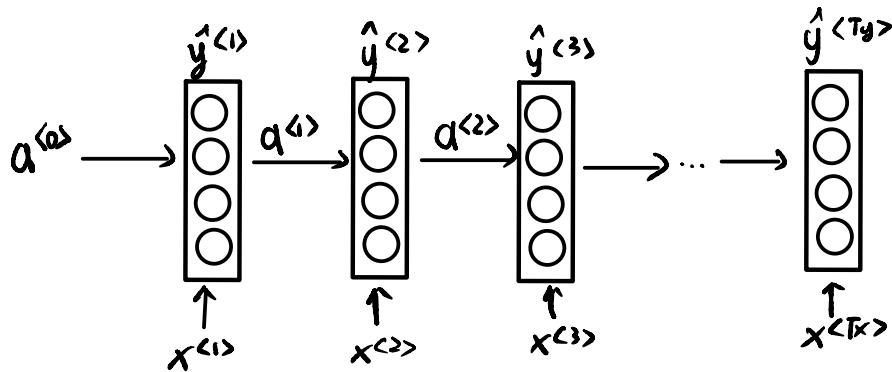
Why standard network doesn't work well with sequence data?



Problems:

- Inputs, outputs can be different lengths in different examples
- Doesn't share features learned across different positions of texts
- Too many weight if input has big shape

RNN



$$\begin{aligned}
 a^{(0)} &= \vec{0} \\
 a^{(1)} &= g_1(W_a a^{(0)} + W_{xa} x^{(1)} + b_a) \leftarrow \text{tanh / sigmoid} \\
 y^{(1)} &= g_2(W_y a^{(1)} + b_y) \leftarrow \text{sigmoid} \\
 a^{(2)} &= g_1(W_a a^{(1)} + W_{xa} x^{(2)} + b_a) \leftarrow \text{tanh / sigmoid} \\
 y^{(2)} &= g_2(W_y a^{(2)} + b_y) \leftarrow \text{sigmoid}
 \end{aligned}$$

Simplified RNN notation

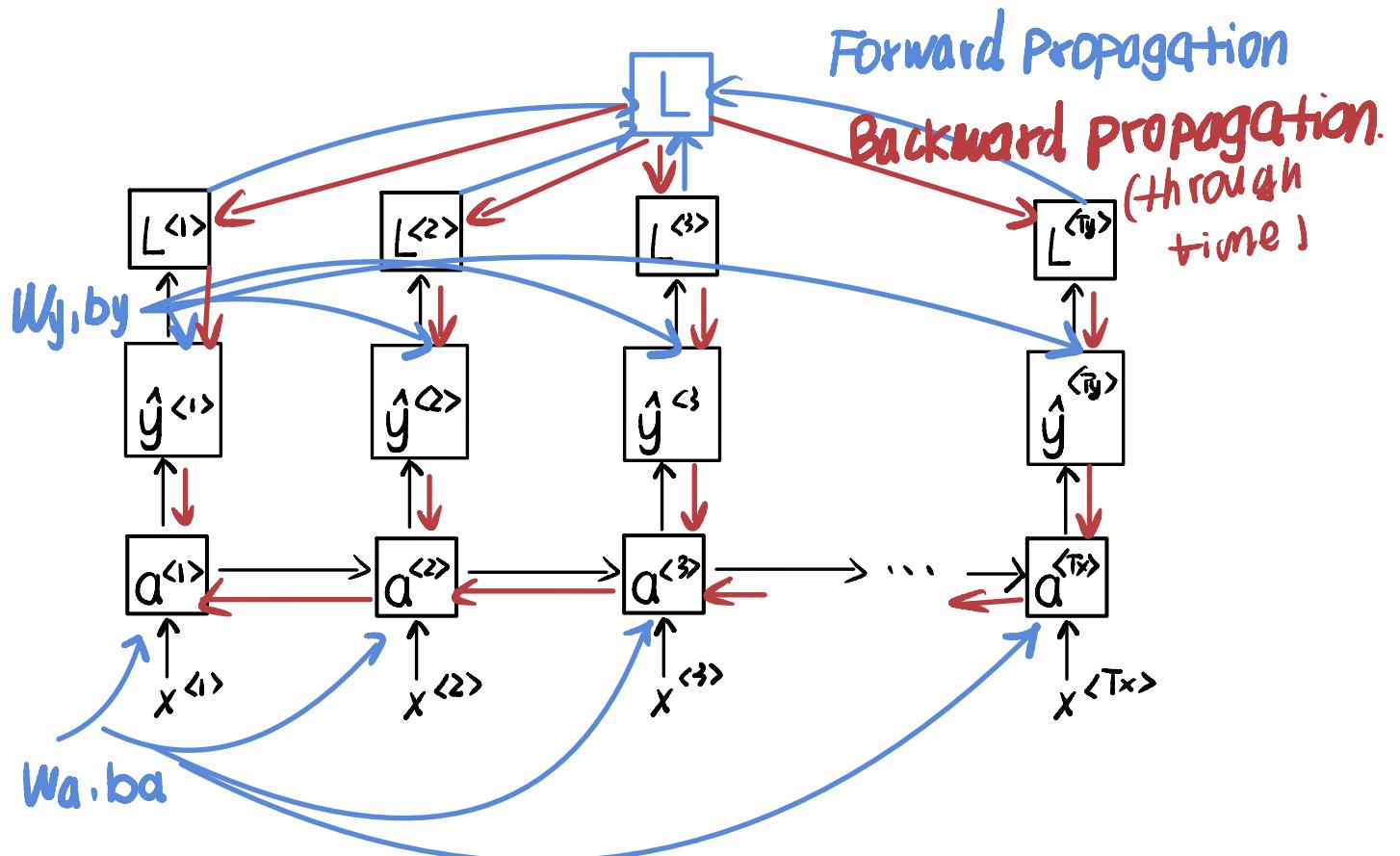
$$\begin{aligned}
 a^{(t)} &= g_1(W_a \underbrace{a^{(t-1)}}_{(100, 100)} + \underbrace{W_{xa} x^{(t)}}_{(100, 1000)} + b_a) \longrightarrow a^{(t)} = g(W_a [\underbrace{a^{(t-1)}, x^{(t)}]}_{\frac{a^{(t-1)}}{x^{(t)}}} + b_a) \\
 y^{(t)} &= g_2(W_y a^{(t)} + b_y) \longrightarrow \hat{y}^{(t)} = g(W_y a^{(t)} + b_y)
 \end{aligned}$$

$\begin{bmatrix} W_a a^{(t-1)} + W_{xa} x^{(t)} + b_a \end{bmatrix} = W_a \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}$

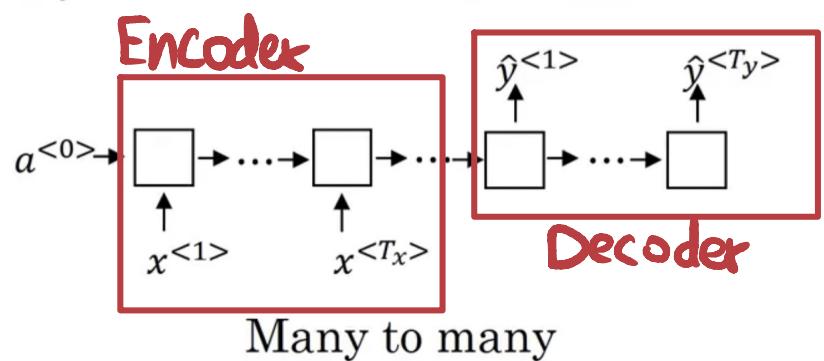
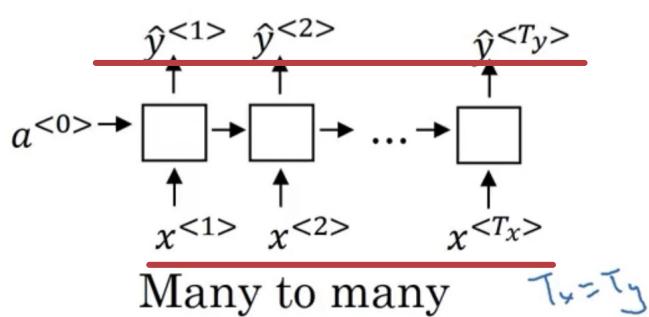
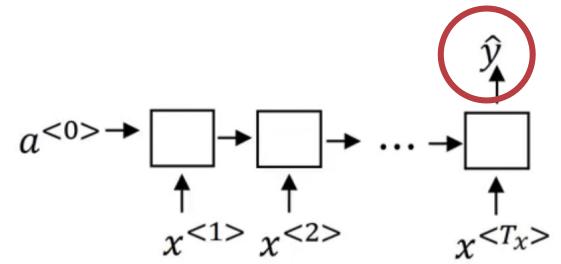
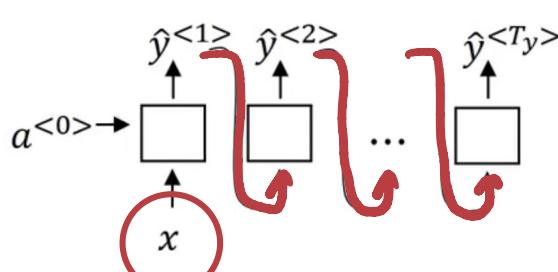
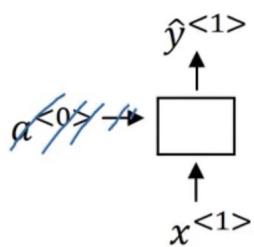
Forward propagation and backpropagation

$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log (1-\hat{y}^{<t>})$$

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$



Different Types of RNN



Language model

Example of language modeling:

The apple and pair salad $\rightarrow P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$

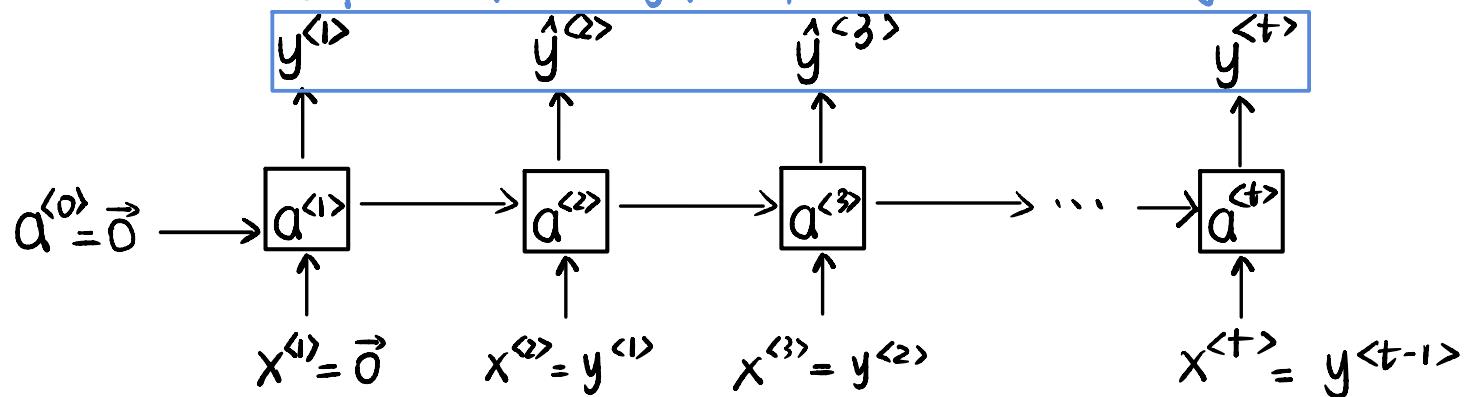
The apple and pear Salad $\rightarrow P(\text{the apple and pear Salad}) = 5.7 \times 10^{-10}$

$P(\text{Sentence}) = ? \quad P(y^{<1>} , y^{<2>} \dots y^{<T>})$

Language Modeling with RNN

Training set: Large corpus of English text

Softmax representing prob of each word in vocabulary

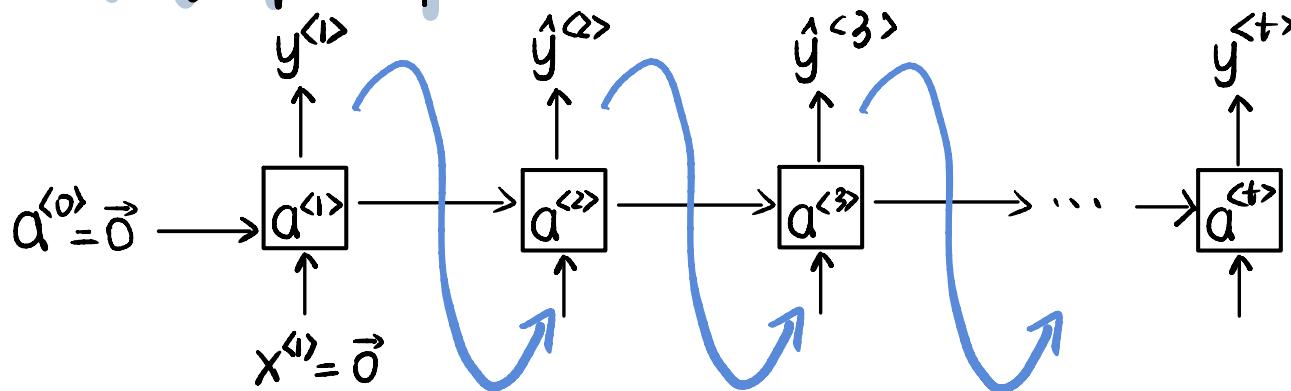


$$l(\hat{y}^{<t>} , y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$L = \sum_t l(\hat{y}^{<t>} , y^{<t>})$$

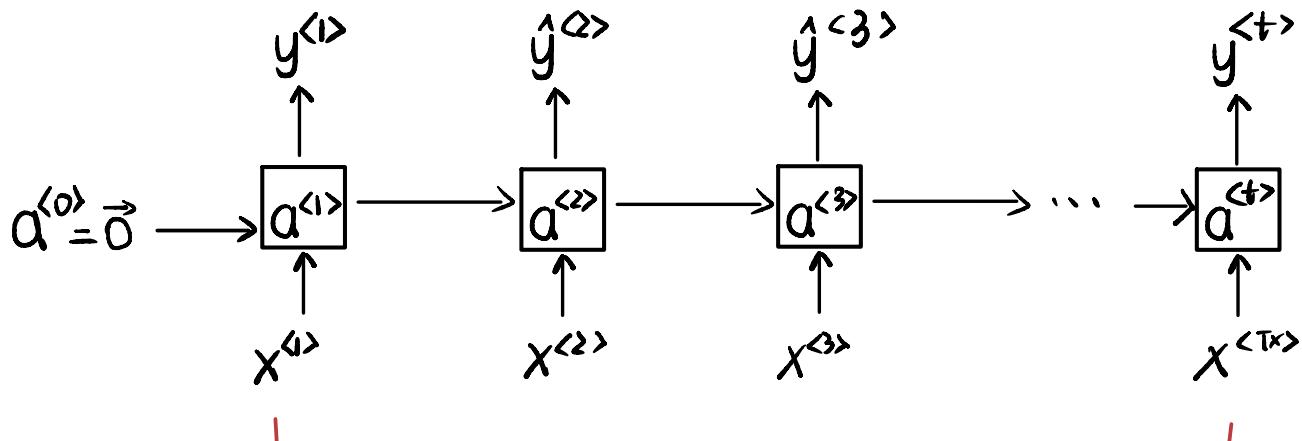
$$\begin{aligned} P(y^{<1>} , y^{<2>} , y^{<3>}) \\ = P(y^{<1>}) P(y^{<2>} | y^{<1>}) \\ P(y^{<3>} | y^{<1>} , y^{<2>}) \end{aligned}$$

Sampling sequence from a trained CNN



GRU And LSTM

Vanishing gradients with RNNs



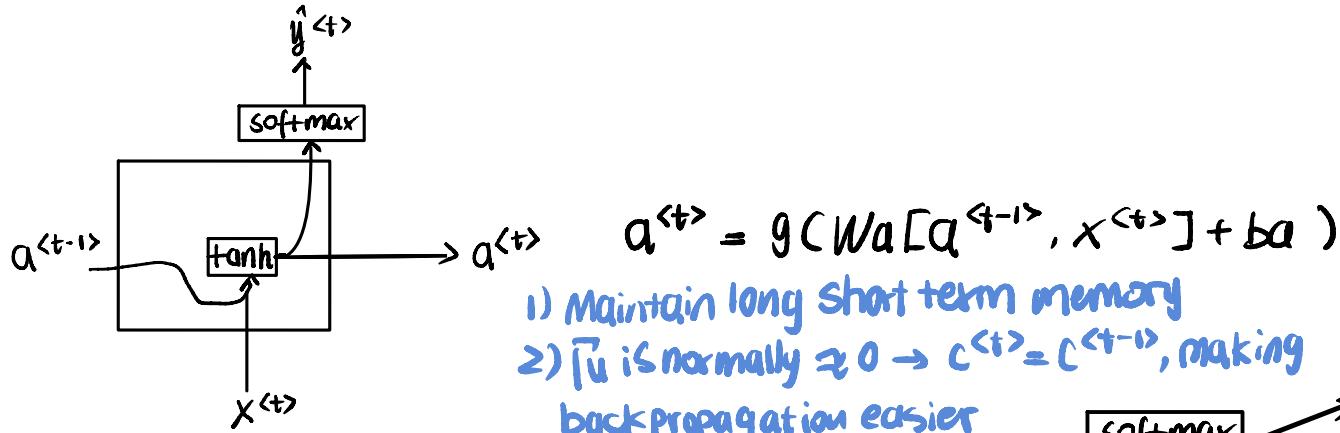
→ Vanishing / exploding gradients.

e.g. The cats which , were full
 The cat which , was full.
 In English, this part can be arbitrarily long.

this relation can be buried
too deep to be discovered
remembered by network

Gated Recurrent Unit (GRU)

Normal RNN unit:



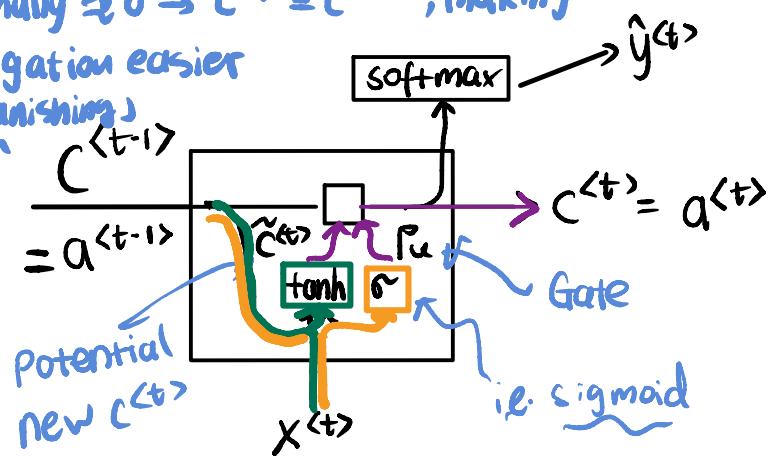
GRU (Simplified): (less gradient vanishing)

$$\tilde{c}^{(t)} = \tanh(W_c[c^{(t-1)}, x^{(t)}] + b_c)$$

$$\tilde{r}_u = \sigma(W_u[c^{(t-1)}, x^{(t)}] + b_u)$$

$$c^{(t)} = \tilde{r}_u * \tilde{c}^{(t)} + (1 - \tilde{r}_u) * c^{(t-1)}$$

↑ element-wise



Full GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

LSTM: long short Term Memory

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

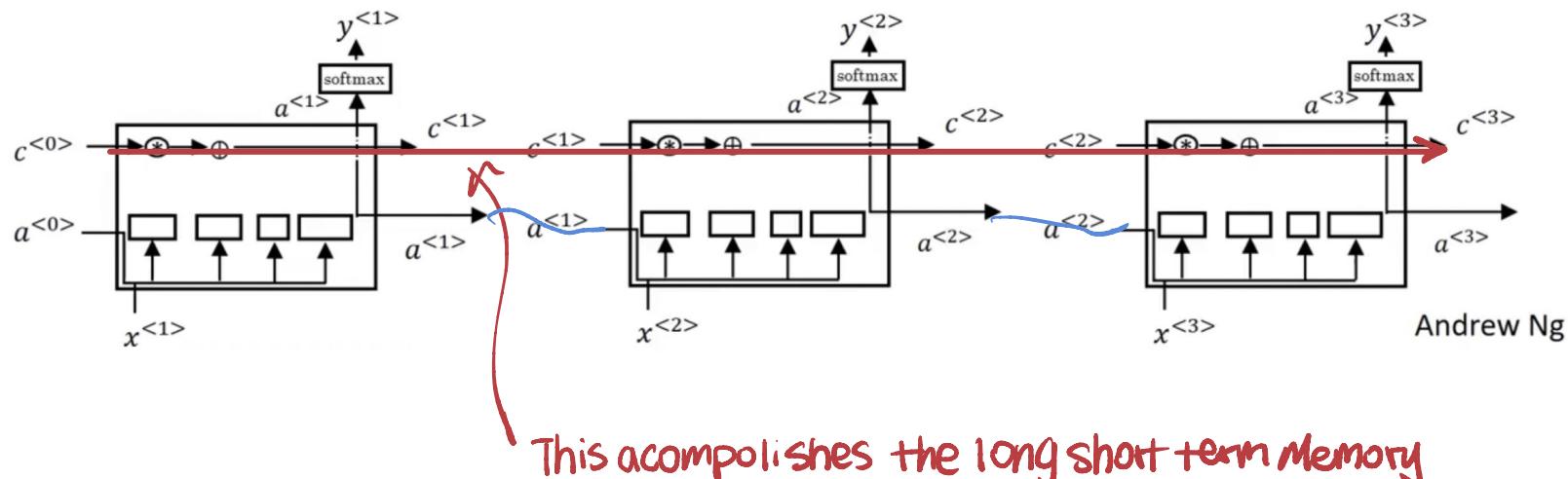
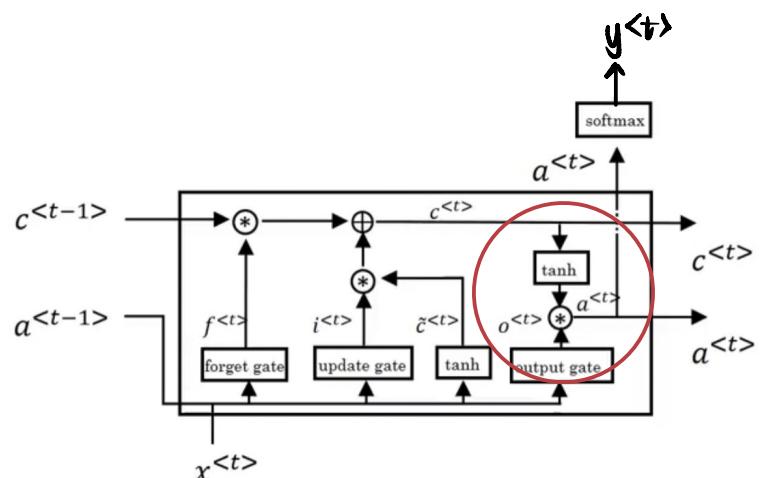
$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$



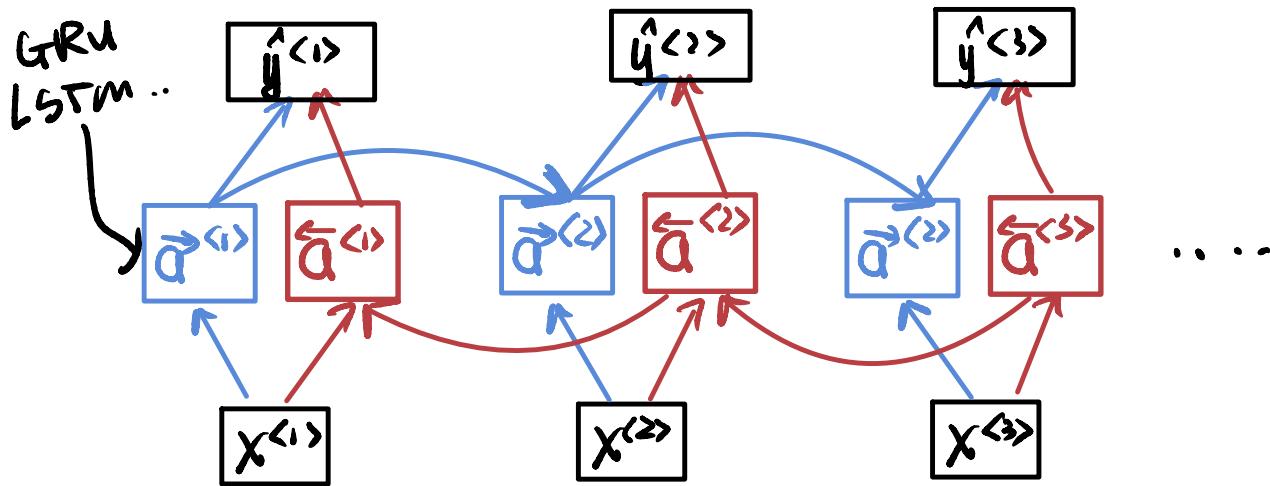
Bidirectional RNN

Motivation:

Unidirectional data is not enough!

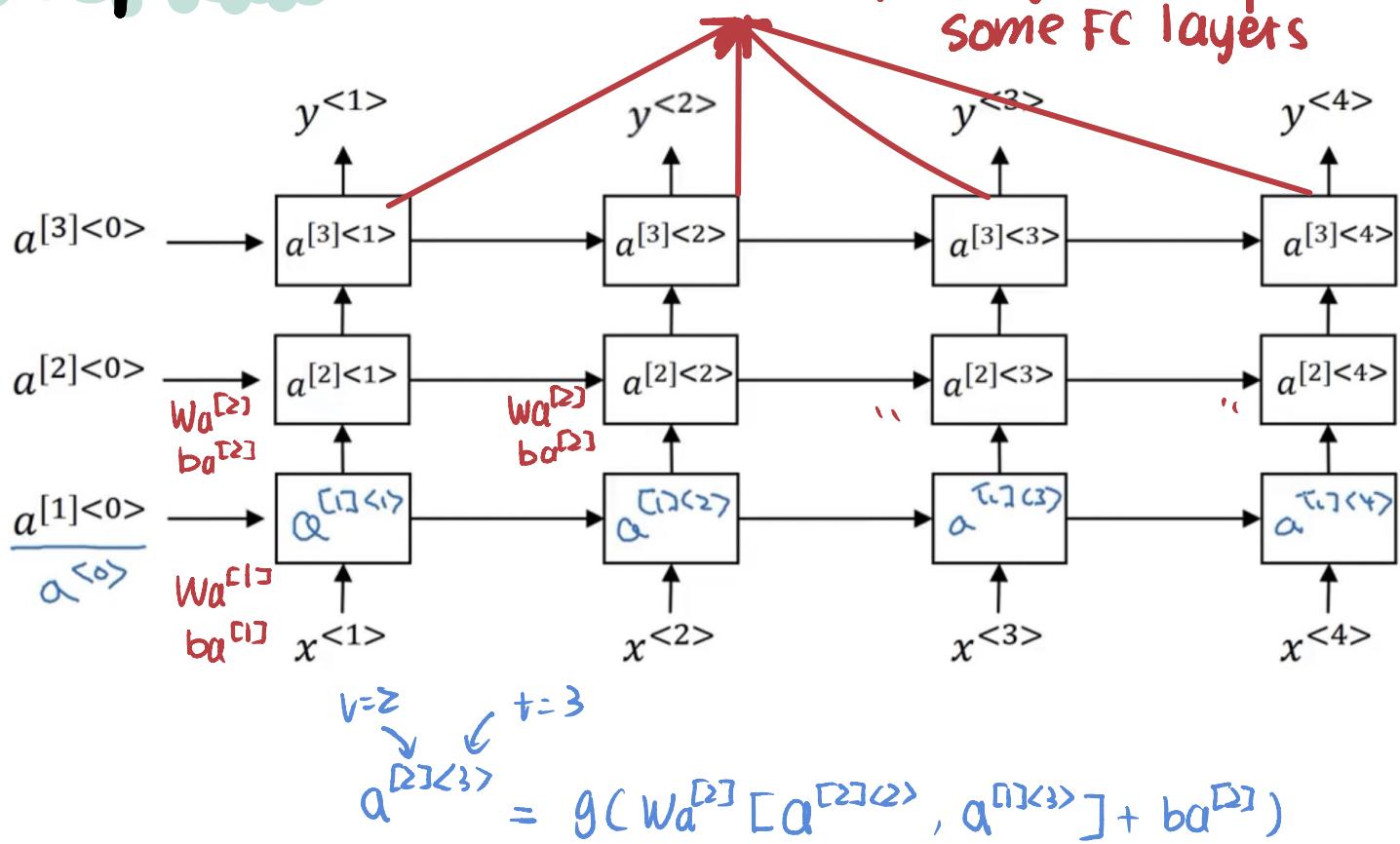
He said, " Teddy bears are on sale ! "

He said, " Teddy Roosevelt was a great President! "



Deep RNN

The last few layers can potentially have some FC layers



Andrew Ng